

Examen psychologique

La psychotechnique des aptitudes. Pour différencier une  
sociotechnique de l'évaluation sans mesurage et une  
psychologie balbutiante de la compréhension de la  
performance

*Mental testing: Differentiating sociotechnical assessment without  
measurement and scientific explanation*

S. Vautier

*Octogone, université de Toulouse, 5, allées Antonio-Machado, 31058 Toulouse cedex 9, France*

Reçu le 21 janvier 2014 ; accepté le 23 janvier 2015

---

**Résumé**

Une conception répandue consiste à considérer que des tests psychotechniques validés permettent de mesurer des aptitudes intellectuelles à partir du scorage des performances observées. Cet article développe une conception falsifiable du mesurage ordinal impliquant que les performances observées falsifient vraisemblablement cette conception et analyse comment la modélisation psychométrique satisfait l'impératif comparatif qui sous-tend l'évaluation des aptitudes. Mais l'efficacité évaluative s'établit au détriment de la connaissance scientifique des déterminants de la performance. La pratique de l'examen psychologique est ensuite analysée comme une sociotechnique de l'évaluation sans mesurage.

© 2015 Société française de psychologie. Publié par Elsevier Masson SAS. Tous droits réservés.

*Mots clés* : Tests psychologiques ; Mesurage ; Psychométrie

**Abstract**

A widespread view consists in considering that validated psychotechnical tests enable one to measure intellectual abilities with the help of the scoring of observed performances. This paper (i) elaborates a falsifiable conception of ordinal measurement, (ii) shows that it is likely that the observed performances falsify it, and (iii) analyzes how psychometric modeling fulfils the comparative imperative that underpins

---

Adresse e-mail : [vautier@univ-tlse2.fr](mailto:vautier@univ-tlse2.fr)

<http://dx.doi.org/10.1016/j.prps.2015.01.005>

1269-1763/© 2015 Société française de psychologie. Publié par Elsevier Masson SAS. Tous droits réservés.

the assessment of abilities. But the evaluative efficacy builds up to the detriment of scientific knowledge of the performance's determinants. The practice of psychological assessment is then thought of as a sociotechnics of assessment without measurement.

© 2015 Société française de psychologie. Published by Elsevier Masson SAS. All rights reserved.

*Keywords:* Psychological testing; Measurement; Psychometrics

*La crainte qu'éprouve le fils authentique de la civilisation moderne à l'idée de s'éloigner des faits qui sont déjà schématiquement préformés par les conventions dominantes de la science, du commerce et de la politique, est la même que la crainte qu'inspire la déviation sociale.*

Max Horkheimer et Theodore W. Adorno, *La dialectique de la raison*.

## 1. Introduction

Tout praticien des tests d'aptitude sait que les scores ou les notes qu'il attribue aux personnes qu'il teste ne mesurent pas de grandeur analogue à la longueur ou la température d'un corps. Il sait aussi l'importance de la locution « test validé » : on n'entreprendrait pas sans risque professionnel une évaluation des aptitudes avec des tests non validés. Cette quasi-labellisation est parfois considérée dans la communauté des utilisateurs comme un gage de scientificité (Gaillard, Colasse, Guihard, & Michel, 2011, p. 155). Une telle opinion est discutable (Lacot, Afzali, & Vautier, à paraître). Si un test validé ne mesure rien, la finalité du scorage psychotechnique doit être assumée dans une perspective sociotechnique par opposition à scientifique (pour la différence entre science et technique, voir Granger, 1995) ; ce qui entraîne des conséquences politiques puisqu'il s'agit de dénaturer l'objet de l'évaluation psychotechnique en lui reconnaissant le caractère d'un fait social par opposition à un fait brut (cf. Searle, 1995).

La position mise à l'épreuve dans cet article est la suivante : l'évaluation d'un niveau de performance ou d'aptitude<sup>1</sup>, à l'aide d'un score qu'on rapporte à une norme statistique, ne constitue pas une opération de mesurage, ce qui implique que l'utilisation du terme de mesure est trompeuse. Si la communauté des utilisateurs de tests souhaite assumer sa responsabilité scientifique, elle doit « faire le ménage » dans ses modes d'expression pour clarifier son domaine de compétences, autant vis-à-vis de ses membres que des membres de la société civile au sens large, en évacuant de sa terminologie les termes connotant la mesurabilité des grandeurs psychologiques, et en assumant l'évaluation comme un processus qui assigne à la personne une ou des propriétés extrinsèques.

Le terme d'évaluation possède une ambiguïté descriptive et appréciative redoutable. On dit qu'on évalue l'intelligence ou la taille d'un enfant. La taille est un attribut mesurable. D'où la tentation de conclure que l'intelligence est mesurable, puisqu'on l'évalue. Or les tests d'aptitude ne sont pas faits pour déterminer une quantité d'intelligence, mais pour assigner une place à l'enfant (via sa performance) dans une échelle de scores.

<sup>1</sup> Le praticien vise l'aptitude en regardant la performance.

En soi, un score n'est pas un jugement de valeur, mais ce n'est pas non plus le résultat d'un mesurage (Vautier, 2014d). Le score, interprété comme une propriété de l'enfant (le score traduit la performance, laquelle est le produit de l'intelligence, donc le score décrit l'intelligence), constitue la condition de possibilité pour que la formulation d'un jugement de valeur sur l'enfant, en fonction du contexte dans lequel il s'agit de l'insérer, acquière une factualité suffisante.

« Insérer une personne dans un contexte social », en l'occurrence l'enfant évalué, implique une construction de significations à propos de la personne qui devient objet d'attention, objet à spécifier, à positionner et vis-à-vis duquel se positionner, objet à insérer dans un réseau d'enjeux relationnels et/ou institutionnels, le plus souvent implicites. Par exemple, Binet et Simon (1907) proposent leur « échelle métrique de l'intelligence » pour « une situation où des doutes planent sur les causes du retard scolaire » (p. 92) et où l'enjeu consiste à « envoyer l'élève à la classe de perfectionnement » ou bien à le renvoyer « à l'école ordinaire ». La performance de l'enfant doit alors être la variable d'une fonction compatible avec l'évaluation, c'est-à-dire d'une fonction dont les valeurs sont compatibles avec les notions pratiques de « pas assez », « trop », « suffisamment » : la performance doit être suffisamment élevée pour un renvoi à l'école ordinaire, ou bien assez basse pour une orientation en classe de perfectionnement (Vautier, 2014b).

Comme la description de la performance n'est pas un nombre (on verra que c'est généralement un *m*-uplet), le scorage prépare son évaluation en transformant sa description en scalaire, ou encore, en nombre, toujours lisible comme degré dans un ordre simple<sup>2</sup>, auquel il suffit d'adjoindre des seuils ajustables à la situation.

Ainsi peut-on exhumer l'impératif de comparabilité de la pratique du scorage psychotechnique. La psychotechnique répond à une demande sociale de comparabilité. Non pas qu'il s'agisse de comparer à tout va ; ce qui importe, c'est de pouvoir comparer si le besoin s'en fait sentir. L'intérêt social de la psychotechnique comme technicité qui s'exerce sur autrui dépend de sa capacité à satisfaire l'impératif de comparabilité. De ce point de vue, on peut saluer la clarté avec laquelle Reuchlin (1969) définit la finalité des tests : « ils fournissent les moyens d'exprimer ces observations [les réponses] sous une forme telle que soient possibles la comparaison [des] individus entre eux et la comparaison de chacun avec les “normes” (descriptives) de la population à laquelle ils appartiennent » (p. 22).

Huteau et Lautrey (1999, p. 76) écrivent que la mesure de l'efficacité intellectuelle – via l'observation de performances à des items de tests – est fondée au niveau ordinal. C'est faux : la psychotechnique du scorage fabrique la comparabilité des performances au lieu de la révéler ; c'est une technique (ou une ingénierie) sociale qui n'exploite aucune loi psychologique connue ni, a fortiori, aucun principe de mesurage.

Cet article définit ce que serait un mesurage ordinal en prenant l'exemple d'un test connu, et montre comment le discours psychométrique entérine le fait qu'on ne sache mesurer ordinalement aucune grandeur théorique avec des réponses (ou des performances) à des items de test. Puis, il analyse les pratiques linguistiques en cours dans la littérature psychotechnique pour montrer comment l'emploi des mots masque ce fait, en prenant comme exemple le manuel du test. Cette analyse est complétée d'une petite mise en scène qui vise à rendre sensibles les tensions logiques et éthiques que doit affronter le psychologue clinicien lorsqu'il sert la démarche évaluative.

<sup>2</sup> Un ordre simple est un ensemble dont les éléments pris par paires peuvent toujours être ordonnés l'un par rapport à l'autre (plus, moins, ou aussi que).

## 2. Mesurer une grandeur avec le test Cubes du WISC-IV

Soient des conditions suffisantes pour le mesurage ordinal d'une grandeur théorique dans une certaine population d'unités d'observation. Ces conditions forment une hypothèse théorique qui est fautive. Par conséquent, l'argument selon lequel on dispose d'une hypothèse de laquelle déduire qu'on sait comparer des quantités d'une grandeur dans cette population est logiquement valide, mais il est faux parce que l'hypothèse sur laquelle il repose est fautive. Il en résulte qu'en l'absence d'hypothèse alternative, on ne sait pas justifier qu'on sache mesurer ordinalement la grandeur théorique avec la performance dans cette population. Cette grandeur n'est pas un concept scientifique.

Le WISC-IV est une batterie de tests utilisée par les psychologues cliniciens dans le cadre de l'examen psychologique de l'enfant et de l'adolescent (Chartier & Loarer, 2008 ; Grégoire, 2009 ; Jumel & Savournin, 2013). Elle comprend 15 tests et permet de calculer, en fonction des réponses observées, des scores, appelés notes ou indices, de « compréhension verbale », de « raisonnement perceptif », de « mémoire de travail », de « vitesse de traitement », ainsi qu'une note « totale » (Wechsler, 2005a). L'analyse porte sur l'hypothèse de mesurabilité d'une grandeur théorique par la performance observée au test Cubes.

### 2.1. La description de la performance au test Cubes

La description de la performance au test Cubes mobilise un langage dont il est utile de connaître la syntaxe et le lexique. Tout d'abord, comme le test comprend 14 tâches, la performance au test est un 14-uplet. Le vocable de *m*-uplet est fondamental pour la compréhension de ce qui suit, c'est pourquoi il convient de s'y attarder quelque peu en partant d'un exemple (voir aussi Vautier, 2014c).

Supposons pour simplifier que le test ne comprenne que trois tâches, toujours administrées dans le même ordre. La performance au test est alors décrite sous la forme d'un triplet (un 3-uplet), par exemple le triplet (1, 1, 0), qu'on peut abrégé par la notation « 110 ». Le premier « 1 » spécifie le résultat issu de l'observation de l'enfant lorsqu'il est confronté à la première tâche ; le deuxième « 1 » spécifie le résultat issu de l'observation de l'enfant lorsqu'il est confronté à la seconde tâche ; le « 0 » spécifie le résultat issu de l'observation de l'enfant lorsqu'il est confronté à la troisième tâche. La syntaxe de la description de la performance à ce petit test prend la forme « 1 puis 1 puis 0 ». Après la syntaxe, penchons-nous sur le lexique de la description. Le résultat de l'observation de l'enfant face à une tâche s'exprime dans un lexique spécifique. Si on considère qu'une tâche est échouée ou réussie, le triplet « 110 » indique deux réussites successives puis un échec grâce au codage « 0 = échec » et « 1 = réussite ».

La description de la performance aux 14 tâches du test Cubes est un 14-uplet. Cette description est multivariée (plus précisément 14-variée)<sup>3</sup>. Le test Cubes comprend trois lexiques. Le premier lexique s'applique pour la description du résultat obtenu à chacune des trois premières tâches. Le second lexique s'applique pour la description du résultat à chacune des tâches n° 4 à 8. Le troisième lexique s'applique pour la description du résultat à chacune des six dernières tâches.

Le premier lexique comporte trois modalités qui sont « 0 » pour « échec », « 1 » pour « réussite partielle » et « 2 » pour « réussite totale ». Ainsi, la performance aux trois premières tâches est

<sup>3</sup> La passation du test Cubes obéit à une règle de départ et une règle d'arrêt, ce qui signifie que dans certaines conditions, la performance n'est pas un 14-uplet, auquel cas le test n'est pas complètement standardisé. Il est inutile de tenir compte de cette particularité ici.

décrite par un triplet parmi les 27 triplets possibles 000, 001, . . . , 222. Le second lexique a deux modalités qui sont « 0 » pour « échec » et « 4 » pour « réussite ». L'utilisation du chiffre 4 au lieu du chiffre 1 pour coder la réussite indique la valorisation de la réussite : réussir une tâche parmi les tâches n° 4 à 8 vaut mieux que réussir une tâche parmi les tâches n° 1 à 3. Le troisième lexique comprend cinq modalités : « 0 » signifie « échec » et les chiffres « 4 », « 5 », « 6 » et « 7 » signifient des degrés croissants de réussite.

Il n'est pas nécessaire pour le propos de préciser comment le psychologue utilise ces lexiques ; supposons seulement que les psychologues qui pourraient effectuer la description d'une certaine performance (qu'on aurait filmée par exemple) soient interchangeables – le test est réputé « cotation-objectif ».

## 2.2. *Le principe de mesurage : éléments théoriques*

Ce qui précède définit le cadre descriptif des phénomènes empiriques qu'on peut décrire avec le test. Voyons maintenant comment on peut imaginer un principe général permettant de relier la grandeur théorique visée par le test à l'ensemble des performances observables. Par « observables », il faut entendre « qui peuvent être observées lorsqu'on procède à une observation », par opposition à l'énumération de toutes les possibilités logiques générées par le langage descriptif de la performance, qui ne dépend d'aucune observation, et qui constitue l'ensemble de réponses logiquement possibles, par opposition à l'ensemble des réponses empiriquement possibles (i.e., celles qui s'observent en fait).

Il faut imaginer un principe pour chaque tâche, avant de résoudre le problème à l'échelle de la description 14-variée. Comme il s'agit de démontrer que la construction théorique requise pour fonder l'idée que la performance au test mesure une grandeur psychologique est fautive, il ne sera pas nécessaire de développer toute la démarche. Il suffit d'en développer une partie et de montrer que cette partie est fautive pour que la théorie complète, qui contient la théorie partielle, soit fautive.

Considérons une tâche dont le résultat est décrit par les chiffres 0 ou 4, qui constituent le lexique le plus parcimonieux du test Cubes. On suppose une grandeur psychologique dont la variation détermine la variation du résultat à la tâche n° 4. Cette grandeur possède par hypothèse une origine naturelle, qu'on peut noter O, en posant qu'elle désigne l'absence de quantité – on admet qu'une quantité négative d'aptitude n'existe pas. On peut aussi considérer que la grandeur possède un maximum qu'on notera « max ».

Le problème consiste à définir une relation du segment [O, max] dans l'ensemble contenant les éléments 0 et 4. À tout point de la grandeur, on veut faire correspondre une réponse observable, 0 ou 4. On veut aussi que tout point de la grandeur ne corresponde qu'à une réponse, sinon cette relation ne pourrait pas être utilisée comme un principe de mesurage. On veut donc une application de [O, max] dans {0, 4}. Enfin, comme la valeur descriptive « 4 » indique par définition un niveau supérieur à celui qu'indique la valeur descriptive « 0 », cette application doit être croissante.

La seule solution possible est une fonction par palier. Ainsi, dire que la réponse à la tâche n° 4 mesure la grandeur revient à invoquer une fonction à deux paliers, les deux paliers étant séparés par un seuil dans [O, max], dont on ignore la position. Lorsque, par une expérience de pensée, on fait varier la grandeur de O jusqu'au seuil, on pose qu'on observe « 0 » ; quand la grandeur dépasse le seuil et varie jusqu'à son maximum, on pose qu'on observe « 4 ». Cette construction théorique n'est pas falsifiable, puisqu'on ne connaît pas la valeur de la grandeur et qu'on peut observer 0 ou 4. Mais elle fournit un cadre logique pour relier intelligiblement le lexique descriptif de la réponse à la tâche et la grandeur que la réponse est supposée mesurer.

On applique la même démarche à la tâche n° 5, en inventant un autre seuil. La question qui se pose maintenant est de savoir comment ordonner les deux seuils sur le segment  $[O, \max]$ , étant donné qu'on suppose que les deux tâches mesurent la même quantité théorique. Notons A et B les deux seuils respectifs. Le langage de la grandeur implique que soit  $A < B$ , soit  $A = B$ , soit  $B < A$ . Comme la tâche n° 4 est supposée plus facile que la tâche n° 5, A se trouve avant B. En effet, l'ordre de difficulté des deux tâches implique la possibilité qu'un enfant possède une quantité d'aptitude telle qu'elle lui permet de réussir la tâche n° 4 mais pas la tâche n° 5. Dans ce cas, cette quantité est supérieure à A et inférieure à B. Donc A est inférieur à B. En d'autres termes, la performance (4, 0) signifie que la quantité d'aptitude de l'enfant se trouve après A – d'où le « 4 » de (4, 0) – et avant B – d'où le « 0 » de (4, 0).

Une conséquence capitale découle de ce qui précède. Cet enfant ne peut théoriquement pas exhiber la performance (0, 4), puisque s'il réussit la tâche n° 5, c'est que sa quantité d'aptitude est supérieure au seuil B, et donc qu'elle est aussi supérieure au seuil A<sup>4</sup>. D'après la fonction par palier de la tâche n° 4, on devrait observer une réussite et non pas un échec à la tâche n° 4.

### 2.3. La falsifiabilité du principe de mesurage et ses conséquences techniques

Nous disposons d'un cadre logique pour relier la grandeur théorique et la performance observable avec deux items, et ce cadre d'interprétation possède un falsificateur, qui est l'observation (0, 4). Donc la théorie est falsifiable, ou encore testable (Popper, 1973). Si on admet que la quantité théorique peut varier lorsque l'enfant passe d'un item à l'autre, la théorie n'est plus falsifiable mais tautologique. Supposons qu'on considère maintenant que la théorie s'applique à tout enfant satisfaisant un certain nombre de conditions (conditions initiales). On peut alors énoncer la loi suivante : quel que soit un enfant dans ces conditions, il ne peut pas produire la performance (0, 4) puisque la performance mesure sa quantité théorique selon la fonction de mesurage que nous venons d'élaborer. Autrement dit, nous venons de dire que la probabilité d'observer l'événement (0, 4) dans ces conditions est nulle (pour une élaboration de la notion de loi en psychologie, voir Vautier, 2011, 2013 ; Vautier, Lacot, & Veldhuis, 2014).

Supposons que des observations, nombreuses, corroborent cette prédiction – aucun « 04 » n'a été observé. Alors le langage de la grandeur théorique est une commodité linguistique pour énoncer cette loi empirique<sup>5</sup> de manière concise (pour une analyse de la fonction descriptive de la théorie en physique, voir Duhem, 2007). Nous ne savons pas si la grandeur existe en tant que telle, nous savons seulement que la fonction de mesurage dont elle constitue le domaine de définition est un modèle commode et prometteur. Supposons maintenant que quelques observations falsifiantes aient été rapportées dans la littérature scientifique<sup>6</sup>. Un nouveau problème scientifique se pose : de quoi d'autre que la quantité théorique dépendent de telles observations ? Quelles que soient les solutions envisageables, l'existence du problème crée un impératif technique : ce qui est observé est une anomalie au regard de la théorie de mesurage ; un argument qui interprète ce qui est observé en termes de niveau de la grandeur théorique n'est pas valide. L'exploitation de la technique de mesurage doit prendre en compte le fait que parfois, les données sont aberrantes. Il ne s'agit pas d'une erreur de mesure au sens d'un manque de précision conduisant à la nécessité d'utiliser

<sup>4</sup> Ici, on a besoin de postuler que la quantité que l'on veut mesurer varie de manière négligeable entre le moment où l'enfant traite l'item n° 4 et le moment où il traite l'item n° 5.

<sup>5</sup> Cette loi structurale est aussi connue sous le nom d'échelle de Guttman (1944).

<sup>6</sup> Ce qui, soit dit en passant, est quasi impossible si les politiques éditoriales des revues d'évaluation en psychologie décèlent que ce type d'étude manque de portée.

un encadrement de la valeur théorique plutôt qu'une valeur ponctuelle, mais d'une aberration théorique qui nécessite une élucidation parce qu'elle signale qu'on ne comprend pas ce qui se passe. Dès lors, une précaution élémentaire consiste à ne pas qualifier ces observations comme des données exploitables et l'utilisateur doit affirmer clairement qu'il ne peut rien conclure de ses observations parce qu'elles sont théoriquement inintelligibles – la performance ne dépend pas que de la quantité théorique, donc la théorie ne « marche pas ». Supposons enfin que de nombreuses observations falsifiantes aient été rapportées. Alors, l'intérêt scientifique de la théorie de mesurage est négatif : on a appris qu'une telle construction théorique est fautive, ce qui constitue une authentique connaissance scientifique.

#### 2.4. *Incertitudes à propos de l'incertitude*

L'incertitude est une notion vague tant qu'on ne précise pas sur quoi elle porte. Lorsqu'on dispose d'un modèle de mesurage ordinal corroboré qui est fondé sur des observations multivariées, on dispose d'une théorie générale dont la vérité est incertaine. La généralité de la théorie est limitée à la population des êtres qu'on peut évaluer. Par exemple, la proposition « quels que soient les enfants qui rempliraient certaines conditions (l'âge, et d'autres attributs descriptifs liés à la manière dont se déroule la passation du test), la performance observée serait 00, 40 ou 44 » est une proposition générale. Comme la proposition est contrefactuelle, le nombre d'unités d'observation est infini et on ne peut donc pas vérifier la proposition unité par unité. On sait au mieux qu'un certain nombre de tests (au sens poppérien du terme) corroborent cette proposition. Face à l'incertitude irréductible de cette proposition, on se contente de considérer qu'elle est vraie jusqu'à preuve du contraire.

Supposons qu'on décide de croire en une telle loi parce qu'elle a toujours été corroborée. Une autre incertitude s'y attache, qui prend la forme d'une indétermination intrinsèque. La gradation des performances 00, 40 et 44 définit trois segments sur  $[0, \max]$ , dont on connaît l'ordre mais pas l'étendue, ce qui implique qu'on est rigoureusement incapable de justifier que le score soit une mesure quantitative (ou additive). Si on additionne les chiffres dans les couples 00, 40 et 44, on obtient respectivement 0, 4 et 8, mais il est évident que la proposition «  $0 + 4 = 4$  », par exemple, est scientifiquement absurde bien que mathématiquement vraie. L'addition n'a pas de sens psychologique. Ces « scores » signifient seulement que la quantité détectée par l'observation « 0 » (i.e., 00) est plus petite que la quantité détectée par l'observation « 4 » (i.e., 04), qui est elle-même plus petite que la quantité détectée par l'observation « 8 » (i.e., 44). Supposons enfin qu'on augmente le nombre d'items jusqu'à un nombre  $m$  et qu'il soit possible d'identifier une fonction par palier à  $m$  seuils (le nombre de seuils dépendant du nombre de valeurs descriptives associées à ces items). On aura affiné le grain de l'échelle ordinale, mais la mesure demeurera ordinale, c'est-à-dire que l'addition des scores demeurera une absurdité psychologique (ou scientifique).

#### 2.5. *L'ambiguïté scientifique de la modélisation psychométrique*

Les psychométriciens savent bien que la fonction par palier qui est nécessaire pour fonder l'idée d'un mesurage ordinal d'une grandeur à l'aide d'une performance multivariée est fautive (Bertrand, El Ahmadi, & Heuchenne, 2008 ; Borsboom, 2008). Il serait tout de même surprenant qu'un phénomène aussi complexe qu'une performance à un test d'aptitude obéisse à un principe si simple, qui revient à expliquer la performance à l'aide d'une seule « variable latente » (ou théorique). Le fait que le modèle soit faux nous apprend (i) que pour expliquer ces phénomènes,

une théorie plus riche est nécessaire et (ii) qu'il est impossible de déduire de l'observation des performances quoique ce soit en termes de niveau de la grandeur.

Mais, au lieu de prendre acte de ces connaissances pour clamer que le jugement évaluatif ne peut être fondé sur nos connaissances scientifiques faute de mesurage, et, éventuellement, pour encourager un programme de recherche ciblé sur les processus de réponse à des items de tests, les psychométriciens ont conservé l'impératif d'une interprétation unidimensionnelle et quantitative de la performance. Pour ce faire, ils ont modifié le modèle théorique développé ci-dessus en introduisant la notion de probabilité d'observer telle réponse conditionnellement à telle valeur numérique de la grandeur, laquelle est définie sur une échelle d'intervalle grâce notamment au postulat de l'existence d'une fonction caractéristique de l'item (cf. Fischer, 1995). Cette transition est explicitée par Bertrand et al. (2008), ce qui permet d'examiner comment ils la justifient.

Suivons les auteurs pas à pas.

« Si la modélisation de la réussite de sujets à des items veut être réaliste, la théorie précédente est trop abrupte et doit être assouplie. Il est exceptionnel qu'une échelle de Guttman, à cause de sa rigidité, s'applique parfaitement aux données expérimentales » (p. 31).

La signification de l'adjectif « réaliste » est ici non pas descriptive, mais pragmatique. Le réalisme invoqué est en fait un appel à la soumission à la demande sociale d'un savoir fondateur de l'évaluation (sinon, pourquoi être réaliste ?). Du point de vue scientifique, l'existence d'anomalies théoriques est reconnue, mais pas leur fréquence, ni leur caractère falsifiant, ni l'invalidité de l'inférence comparative, « telle quantité est supérieure à telle autre », à partir du modèle et des données. Les auteurs poursuivent en adoptant la position suivante :

« Il est naturel d'interpréter les écarts au modèle en concédant un caractère aléatoire à la relation empirique 'réussir' de S [ensemble des sujets, unités d'observation] vers I [ensemble des items]. Ce caractère aléatoire est dû aux autres variables – non explicitement prises en compte comme la compétence des sujets et la difficulté des items – qui peuvent influencer la réussite ou l'échec ; on imagine aisément qu'elles sont complexes et nombreuses : humeur du sujet, environnement physique et social, mode de présentation de l'item, etc. [...] » (p. 31).

Les auteurs reconnaissent explicitement que les performances dépendent d'une multitude de causes inconnues tant d'un point de vue théorique que pratique. Mais on ne voit pas en quoi cette ignorance implique que la performance observée doive être conçue comme le résultat d'une expérience aléatoire (pour une introduction à la notion d'expérience aléatoire, voir Falmagne, 2003). Il semble que les auteurs confondent les probabilités subjectives, qui servent à jauger la confiance qu'on a en certaines propositions, et les probabilités objectives, qui supposent l'indétermination intrinsèque des phénomènes (Hacking, 2002). Avant d'en tirer les conséquences, poursuivons encore avec eux.

« On conçoit donc que dans la situation où un sujet  $s$  est confronté à un item  $i$ , la réussite de  $i$  par  $s$ , au lieu d'être toujours réalisée quand  $\gamma(s) \geq \delta(i)$ , jamais quand  $\gamma(s) < \delta(i)$  [ $\gamma(s)$  et  $\delta(i)$  désignent respectivement la position de  $s$  et de  $i$  sur le domaine de la grandeur], est gouvernée par une tendance floue : la réussite (et son contraire l'échec) a une certaine probabilité de survenir. Désormais ce n'est plus le fait de réussir, mais les chances de réussir qui seront fonctions de  $\gamma(s)$  et  $\delta(i)$ . La probabilité que  $s$  réussisse  $i$ , notée  $\pi(s,i)$ , dépend de la compétence de  $s$  comme de la difficulté de  $i$ . » (p. 31).



Les anomalies théoriques sont maintenant éliminées par un récit probabiliste qui invente une « tendance floue » à produire telle ou telle performance. Le « flou » s'exprime par des probabilités et ce qu'il y a de permanent dans la « tendance » est sous-tendu par la grandeur nommée « compétence ».

Les psychométriciens ont poussé la rhétorique jusqu'à appeler les modèles psychométriques des « modèles de mesure ». L'intuition quantitative est sauvegardée mais c'est au prix d'un renoncement à la connaissance théorique. La préférence pour la grandeur comme cadre conceptuel assimilateur dépasse l'intérêt pour la compréhension proprement scientifique de la performance, laquelle n'a finalement qu'un rôle auxiliaire. On pourra désormais estimer une valeur numérique, ce qui n'est pas mesurer, quand bien même on admet que cette valeur ne permet pas de comprendre comment la performance a été produite puisque, étant donnée n'importe quelle valeur de la grandeur, toute performance peut être observée—avec une probabilité plus ou moins importante, mais jamais égale à 0 ni à 1. Autrement dit, on accepte l'opacité de la performance et on assigne des valeurs numériques à des performances en toute ignorance<sup>7</sup>. En dépit de « l'évidence », la psychométrie moderne a sauvé l'entreprise comparative du fait qu'on ne sache mesurer de manière ordinaire aucune grandeur psychologique, en substituant l'estimation statistique au mesurage expérimental.

### 3. Les praticiens des tests peuvent-ils revendiquer une fonction d'évaluation et une responsabilité scientifique ?

La partie précédente a montré (i) comment une fonction par palier reliant la grandeur théorique à une performance multivariée permet de comparer des performances distinctes, et (ii) pourquoi un tel modèle est certainement faux, ce qui implique que l'interprétation ordinaire de la performance observée n'est pas acceptable du point de vue logique faute de modèle corroboré. De plus, le recours aux probabilités pour sauvegarder l'intuition quantitative via la modélisation psychométrique voile notre ignorance des déterminants de la performance pour satisfaire un impératif non scientifique qui va être discuté ici (dans un contexte plus large, voir aussi [Pestre, 2013](#), chapitre 3). Ce type d'analyse suggère que l'évaluation psychotechnique des aptitudes constitue un métier extraordinairement ingrat, parce que le praticien doit spécifier comment il articule le besoin social d'assimiler les personnes à des organismes dotés de diverses formes de capacités intellectuelles, conçues comme des grandeurs empiriquement indéterminées mais essentielles pour l'évaluation des individus — « les construits » —, et les questions, toujours ouvertes, (i) de ce qui, dans des conditions particulières (décrites au mieux grossièrement), détermine les performances aux items des tests d'aptitude<sup>8</sup>, et (ii) de ce que ces performances déterminent à leur tour.

L'impératif qui motive la méthodologie du scorage des performances intellectuelles est le même que celui qui institue la notation scolaire, ou qui conduit [Duhem \(2007, partie 2, chapitre 1\)](#), dans une analyse lumineuse de la quantité et de la qualité, à tenter, sans y parvenir, de justifier

<sup>7</sup> En particulier, l'estimation de la valeur numérique de la grandeur à un instant  $t$  pour une personne donnée est logiquement différente de la tendance centrale de la grandeur pour cette personne, si tant est qu'une telle notion admette une interprétation psychologique. Rien n'exclut que la tendance centrale soit 'significativement' différente de l'estimation issue d'une performance ponctuelle. Mais cette incertitude est en quelque sorte renvoyée dans le « pré-conscient épistémologique » du chercheur puisque pour pouvoir l'analyser empiriquement, il faudrait savoir mesurer la grandeur ([Vautier, Veldhuis, Lacot, & Matton, 2012](#)).

<sup>8</sup> « [...] il est peu réaliste de penser qu'on est aujourd'hui capable d'expliciter les mécanismes psychologiques susceptibles de générer les réponses aux items d'un test ou d'un questionnaire » ([Juhel et al., 2011](#), pp. 186–187).

le mesurage de la qualité « être un bon géomètre ». Cet impératif est formulé de façon concise par Perron dans une des discussions de la Conférence de consensus sur l'examen psychologique de l'enfant et de l'adolescent (Voyazopoulos, Vannetzel, & Eynard, 2011) : « [les scores] sont des jugements comparatifs de valeur qui se répercutent au niveau sociologique général, dans une société qui a besoin de hiérarchiser, à l'école ou dans l'entreprise, au niveau de la micro-sociologie et au niveau des jugements de valeur que l'individu porte sur lui-même » (p. 234). [Michell \(2003\)](#) propose une analyse historique de ce qu'il appelle l'impératif quantitatif, mais sans développer la fonction sociale de l'évaluation, laquelle nécessite seulement la projection des performances à évaluer sur une échelle ordinale. [Coombs \(1964, chapitre 13\)](#) formule la nécessité sociale de compresser les observations sur une « ligne de décision ». Cette nécessité est comparative. Si deux performances sont incomparables, alors deux personnes représentées par ces performances sont aussi incomparables, ce qui est rédhitoire pour la demande sociale.

La dynamique de l'examen psychologique des aptitudes repose sur deux aspirations téléologiquement distinctes, évaluer vs comprendre, potentiellement incompatibles. Comme les performances observables ne sont pas simplement ordonnées, les méthodologies déployées sont incompatibles dès lors que la première opère un forçage descriptif par le scorage. On ne peut alors pas évaluer et comprendre la performance dans le même mouvement. Pour l'évaluer, il faut la scorer, c'est-à-dire la faire littéralement disparaître sous le nombre, lequel tirera sa signification d'une référentialisation que [Danziger \(1990, 1987\)](#) qualifie de galtonienne – l'étalonnage du score, ou encore le rapport à une distribution de référence. Tandis que pour comprendre la performance, il faut, adoptant une posture expérimentale, en découvrir les tenants – variables indépendantes – ce qui suppose de s'appuyer pleinement sur le langage descriptif qui permet d'identifier les changements intervenant au niveau de la variable dépendante. Du point de vue temporel, le calcul du score et son interprétation normative (ou, de manière synonyme, évaluative) prennent un instant, tandis que l'investigation de ce qui et de ce que détermine la performance est un effort de pensée, qui conduit peut-être à quelques spéculations ou hypothèses dans le cadre même de l'examen, ou davantage si on dispose de connaissances générales pertinentes – sinon, pas plus, parce que les personnes qui viennent se faire évaluer n'entrent pas de ce fait dans un programme de recherche particulier.

Le fait que la communauté des praticiens de l'évaluation psychotechnique cherche une légitimation de cette pratique dans la doctrine de la validation des tests (e.g., [Juhel et al., 2011](#)) constitue ainsi un obstacle épistémologique ([Bachelard, 1983](#)) au progrès de la connaissance scientifique de ce qui détermine la réponse aux items de tests – on pourrait utiliser la notion de dissonance cognitive. Tout se passe comme si, pour éviter la dissonance, la doctrine de la validation des tests s'était organisée pour ne pas accuser réception de la non mesurabilité des grandeurs psychologiques par des réponses à des tests ; la mesurabilité des grandeurs psychologiques tient le rôle de postulat fondateur, et la signification des notions de mesurage<sup>9</sup>, de validité<sup>10</sup> et de généralité<sup>11</sup>

<sup>9</sup> Mesurer c'est attribuer un nombre.

<sup>10</sup> Un test valide est un test qui mesure bien ce qu'il est censé mesurer ; un argument valide est un argument approximativement vrai ; pour des discussions 'orthodoxes' de la validité en psychologie, voir par exemple [Cizek \(2012\)](#), [Kane \(2006\)](#) ou [Newton \(2012\)](#) et pour des discussions critiques, voir [Borsboom, Cramer, Kievit, Scholten et Franic \(2009\)](#) et [Michell \(2009, 2013\)](#).

<sup>11</sup> Dans la psychologie dite, abusivement, « nomothétique », le général n'est plus ce qui s'applique à toute unité d'une classe de référence, mais ce qui particularise la classe de référence ([Danziger, 1987, 1990](#); [Lamiell, 1998](#); [Salvatore & Valsiner, 2010](#); [Vautier, 2011, 2013](#)).

est adaptée pour les usagers d'un flot sociotechnique, spécialisés dans l'évaluation des aptitudes (Lacot et al., à paraître).

Cette doctrine se focalise sur des propositions qui portent sur les conséquences de l'application du postulat de mesurabilité des grandeurs psychologiques à l'échelle d'agrégats de répondants. L'immense mobilisation académique dédiée à la validation des tests aboutit à proposer des tâches typiques (parfois confidentielles comme dans le cas du WISC-IV), des règles de scoring, des espaces dimensionnels (ou factoriels) associés à des termes parlants mais vagues (les construits – e.g., le Raisonnement Perceptif), et des relations entre variables statistiques. Ce discours est capable d'assimiler toute unité d'observation – toute personne à telle date – comme point parmi une infinité de points possibles, mais la trajectoire individuelle de ces points demeure totalement indéterminée. D'où, dans le contexte de l'examen psychologique qui est une situation individualisée, l'impression que « quelque chose cloche » (Vautier, 2012).

La lecture du manuel d'interprétation du WISC-IV (Wechsler, 2005b) permet de cueillir quelques spécimens du brouillage pratico-conceptuel qui résulte de la subordination réciproque des projets sociotechniques (observer pour évaluer) et scientifique (observer pour comprendre) qui est opérée par la psychotechnique contemporaine, dont le dernier sort perdant parce qu'il devient quasi impossible de renoncer à la mythologie des grandeurs psychologiques que suggère le langage profane.

### 3.1. *Quand le score vaut pour le fonctionnement cognitif*

On lit dans le manuel que « l'usage des notes standard normalisées en fonction des âges permet au praticien de comparer le fonctionnement cognitif de chaque enfant avec celui des enfants du même âge » (Wechsler, 2005b, p. 85). Le fonctionnement cognitif n'est pas une grandeur mais un processus. On ne peut être d'accord avec l'affirmation du manuel sur une base logique. On peut toujours comparer un nombre à un nombre moyen. Par exemple, une note qui se trouve à 1,2 écart type de la moyenne de référence est supérieure à cette moyenne.

Ce fait algébrique implique que l'enfant à qui on vient d'attribuer cette note ait un fonctionnement cognitif supérieur à celui de l'enfant moyen, à condition d'admettre que le score objective une valeur. On peut être d'accord avec le manuel si on prétend évaluer – et non pas mesurer – le fonctionnement cognitif de l'enfant. Dans ce cas, le pouvoir du verbe, qui permet qu'un nombre vaille pour un processus pourvu qu'on se focalise sur la valeur de ce processus, joue à plein ; nous avons quitté le registre de la pensée scientifique parce qu'il n'est plus question de décrire mais d'évaluer.

### 3.2. *Quand le score vaut pour la performance*

Dans la même page du manuel, on lit « une note standard représente la performance d'un enfant à un subtest comparativement à la performance de ses pairs du même âge » (Wechsler, 2005b, p. 85). Le même processus de projection de significations sur le nombre est proposé, puisque cette fois-ci le nombre représente une performance. Nous avons vu que la performance se décrit comme un vecteur (un *m*-uplet), ce qui implique qu'un nombre représente de manière univoque une performance seulement si les performances qu'on peut observer sont simplement ordonnées. Si, comme c'est probablement le cas, les performances qu'on peut observer ne sont pas simplement ordonnées, une note standard représente une performance pour une communauté

d'utilisateurs qui veut bien se donner une telle convention, et alors ce n'est plus la performance qui est représentée mais ce qu'elle vaut dans une échelle de valeurs.

### 3.3. *Le score serait l'apparence trompeuse d'une vérité cachée*

Les promoteurs du WISC-IV écrivent encore que « la note vraie est le reflet de la véritable aptitude du sujet, combinée avec un certain degré d'erreur de mesure » (Wechsler, 2005b, p. 87), en référence à la théorie classique des tests qui postule que le score observé est la somme d'un score vrai et d'une erreur de mesure. Supposons qu'on admette l'existence de la « véritable aptitude du sujet », et qu'on admette aussi que la note vraie, c'est-à-dire la moyenne de tous les scores qu'aurait pu avoir le sujet, mesure – et non pas « reflète » – sa véritable aptitude (laquelle, donc, serait une quantité objective).

Vautier et al. (2014) ont analysé le rôle interprétatif de la notion d'erreur de mesure dans le paradigme néo-galtonien. Dans un modèle psychométrique, l'erreur de mesure garantit qu'on ne pourra jamais encadrer la valeur de la grandeur théorique qu'on projette sur le score de manière falsifiable. Ainsi, l'interprétation psychométrique protège de toute falsification le postulat de l'existence de la grandeur (voir aussi Vautier, 2014a). Ce postulat est intimement lié à la nécessité de l'évaluation, puisqu'il garantit la comparabilité au prix d'une incertitude irréductible due à l'erreur de mesure. La passion d'évaluer, qui projette une grandeur là où on ne sait pas la mesurer, est confortée dans sa toute-puissance par l'acceptation « réaliste » (cf. Bertrand et al., 2008) de la fatalité de l'erreur de mesure inhérente au scorage de la grandeur.

### 3.4. *Le complément d'objet direct du verbe « mesurer » doit être une grandeur*

Le verbe « évaluer » admet des compléments d'objet direct (COD) qui ne sont pas nécessairement des grandeurs, tandis que le verbe « mesurer » n'admet pour COD que des grandeurs. Mais l'utilisateur des tests est acculturé à un discours savant qui emploie les verbes « mesurer » et « évaluer » de façon interchangeable. Le manuel du WISC-IV ne fait pas exception. Par exemple, « l'indice de compréhension verbale du WISC-IV est une mesure de la formation de concepts verbaux, du raisonnement verbal et des connaissances acquises dans le propre environnement du sujet » (Wechsler, 2005b, p. 89). La formation de concepts verbaux etc. ne dénote pas une grandeur. Dans la phrase « l'ICV actuel peut être considéré comme une mesure affinée et plus pure du raisonnement verbal et de la conceptualisation que [...] » (Wechsler, 2005b, p. 89), le raisonnement verbal n'est pas une grandeur, pas plus que le raisonnement perceptif et fluide qu'on trouve dans la citation : « l'indice de raisonnement perceptif est une mesure du raisonnement perceptif et fluide, du traitement spatial et de l'intégration visuomotrice » (Wechsler, 2005b, p. 89)<sup>12</sup>.

## 4. Pourquoi ne pas assumer que l'évaluation psychologique soit une sociotechnique ?

Pour que les praticiens des tests ne soient pas empêtrés dans l'ambiguïté du « désir psychométrique », ils peuvent considérer que la légitimation de l'évaluation psychotechnique relève

<sup>12</sup> On peut aussi ajouter les deux exemples suivants. (1) « l'indice de mémoire de travail procure une mesure de la capacité de la mémoire de travail de l'enfant » (Wechsler, 2005b, p. 90). (2) « l'IVT [indice de vitesse de traitement] fournit une mesure de l'aptitude de l'enfant à inspecter rapidement et correctement des informations visuelles simples, à les traiter de manière séquentielle et à les discriminer [...]. Cette note composite est également une mesure de mémoire visuelle à court terme, d'attention et de coordination visuomotrice » (Wechsler, 2005b, p. 90).

d'une science des contraintes sociales. Par exemple, c'est une contrainte sociale : on n'entre pas dans une formation d'élève pilote de ligne à l'École nationale de l'aviation civile si on occupe une position jugée rédhibitoire dans un certain espace évaluatif. Une telle anthropologisation de l'évaluation psychotechnique, qui se présenterait alors comme une sociotechnique sans mesurage, permettrait au chercheur en psychologie d'étudier la performance aux tests, héritage précieux pour l'observation des performances humaines s'il en est, avec la liberté idéologique d'en interroger tant l'indétermination que la détermination, ce qui supprimerait l'utilité du refoulement des faits falsifiants évoqués plus haut au sein même de la communauté concernée. Elle permettrait au praticien de tourner son attention vers les formes de manifestation du pouvoir social qui s'exercent sur l'individu qu'il inscrit dans un espace de valeurs dont il a l'expertise, et sur lui-même lorsqu'il se fait l'agent de l'évaluation (cf. [Canguilhem, 1958](#)).

Aventurons-nous, par le jeu d'une simulation théâtrale, dans la problématique qui consiste, pour le praticien, à ne pas tricher avec l'enfant qui lui est amené en ce qui concerne ce dont celui-ci est l'enjeu. Un psychologue essaie de dire la vérité en répondant à un enfant curieux. En guise de précaution, ajoutons que ce dialogue n'a pas de vocation normative ni réaliste. Pour les besoins de l'analyse, le psychologue aura la possibilité de consulter la documentation technique du test en présence de l'enfant.

Le test Cubes a été administré à un enfant de huit ans et un mois et la performance observée est la suivante : (2, 2, 2, 0, 4, 4, 0, 0, 0, 0, 0, 0, 0).

L'enfant :

C'est bien ?

Le psychologue :

Ta note à ce test est de 14 points, ce qui fait une note standard de 6 points, d'après ce tableau [il montre la table d'étalonnage qui se trouve p. 210 dans le *Manuel d'administration et de cotation* ([Wechsler, 2005a](#))].

Mais est-ce que c'est bien ?

La plupart des enfants de ton âge réussissent mieux.

Alors c'est pas bien.

Ta note est moins bonne que la note qui sert de point de repère.

C'est parce que je suis pas assez intelligent ?

Qu'est-ce que ça veut dire, être assez ou pas assez intelligent ? Tout dépend de ce que tu veux faire.

Je ne veux pas aller à l'école.

Ce que je peux te dire, c'est que les gens qui ont fabriqué le test l'on fait passer à des enfants qui avaient à peu près le même âge que toi, qu'ils ont calculé la moyenne de leurs notes, et que cette note moyenne est plus grande que ta note.

Alors c'est normal.

Qu'est-ce qui est normal ?

Que ma note ne soit pas la moyenne, parce que tous les enfants ne peuvent pas avoir la même note.

Oui. Ta note est plus petite.

C'est quoi la moyenne ?

C'est 10.

C'est qui ces enfants ?

Je ne sais pas exactement. Les gens qui ont fait le test indiquent, dans ce document [*Manuel d'interprétation*, p. 23 (Wechsler, 2005b)], qu'ils ont fait passer le test à 23 garçons et 23 filles de ton âge.

C'était quand ?

En 2004.

Il y a longtemps. Peut-être que je suis aussi intelligent que 46 enfants de mon âge maintenant. Les gens qui ont fait ce test pensent sans doute que si on faisait passer le test à 46 enfants aujourd'hui, la note moyenne ne changerait pas beaucoup.

Et toi, qu'est-ce que tu en penses ?

Je n'en sais rien. Il y a beaucoup de choses qu'on ignore. Par exemple, on ignore ce qu'est l'intelligence et on ne sait certainement pas la mesurer.

Oui mais tes tests, ils permettent de savoir que je suis moins intelligent que les autres.

Ils permettent de savoir que ta note est plus petite que la note moyenne d'un certain groupe d'enfants. On utilise ce groupe d'enfants comme une image de l'enfant typique de ton âge. J'aime pas passer pour un imbécile et ma note dit que je suis un imbécile à ce test.

Dans la vie, nous sommes tous, à un moment ou à un autre, comparés aux autres. Ta note sert à te comparer à d'autres enfants lorsque tu essaies de résoudre les problèmes du test.

Ah, alors c'est comme sur un podium.

Oui, tu n'es pas parmi les premiers.

Comment tu le sais ?

Parce que j'utilise un modèle qui me dit comment les enfants sont répartis sur le podium.

Mais ce modèle, il est vrai ?

Non. Mais puisque je veux te comparer aux autres, je l'utilise. Les psychologues qui font passer ce test font comme ça, nous utilisons le même modèle comme ça on parle de la même chose.

Même s'il est faux ?

Comme beaucoup de monde considère qu'il n'est pas trop faux, il devient vrai entre nous. Mais toi tu sais qu'il est faux.

Oui.

Alors ce n'est pas grave.

Qu'est-ce qui n'est pas grave ?

Que cette note me fasse passer pour un nul à ce test.

La mise en nombre de la performance au test – le 14-uplet – n'est pas permise par le fait qu'elle serait surdéterminée par une loi d'ordre simple (une loi de structure). Elle constitue cependant la condition de possibilité d'une comparaison sociale de l'enfant lorsqu'il se comporte dans une certaine scène sociale – le test. Si on veut éviter de contredire Huteau et Lautrey (1999) lorsqu'ils affirment que la mesure de l'efficacité intellectuelle est fondée au niveau ordinal, on doit ajouter que ce fondement ne réside pas dans une loi expérimentale mais dans un impératif social : il faut comparer et pour comparer il faut ordonner, d'où les conventions nécessaires – le barème de notation et la normalisation des notes dites brutes.

La conséquence qui découle de cette analyse est que si la demande sociale d'examen psychologique n'exige pas une comparaison sociale, le psychologue n'a pas besoin d'utiliser de notes standard, lesquelles sont la seule justification des notes brutes (qui ne sont pas ontologiquement brutes mais bien socialement construites pour satisfaire l'impératif de comparabilité sociale).

Poursuivons le dialogue entre le psychologue et l'enfant pour illustrer en quoi la demande sociale est une demande de comparaison sociale.

L'enfant :

Tu vas dire mes notes ?

Les gens qui se préoccupent de ton avenir ont besoin que j'aie cette information, mais je ne suis pas obligé de donner le détail.

Alors tant mieux s'ils croient que je suis un imbécile, j'irai pas à l'école.

Ils ont besoin de décider dans quelle école tu vas aller.

Ils vont me mettre dans une école pour imbéciles ?

Je ne sais pas. Ils essaient de trouver quelle est la meilleure solution pour toi parce que ta maîtresse pense que tu risques d'avoir des difficultés si tu restes dans ton école.

J'aime pas cette école.

Penses-tu que tu serais capable d'avoir de bonnes notes si tu restais dans cette école ?

Je sais pas.

Si tu as de bonnes performances aux tests, cela veut dire que tu es capable d'avoir de bonnes notes à l'école et on te fait confiance.

Pourquoi ?

Parce qu'on se dit que, grosso modo, c'est la même intelligence qui te permet de faire les exercices aux tests et de faire les exercices à l'école.

Oui mais t'as bien vu que je suis nul à ce test.

Si tes performances aux tests ne sont pas bonnes, ça ne veut pas forcément dire que tu n'es pas capable d'avoir de bonnes notes à l'école.

J'aimerais bien mais j'y arrive pas.

Qu'est-ce qui bloque la performance de cet enfant ? Quelles contraintes opèrent dans son fonctionnement psychologique ? Ces contraintes permettent-elles de prévoir un échec scolaire dans des conditions de scolarisation dont on ignore les détails ? Répondre à de telles questions en s'appuyant sur des connaissances scientifiques générales, au sens nomothétique du terme (Lamiell, 1998), paraît aujourd'hui prématuré.

Si on se réfère à la modélisation psychométrique, la réponse observée à un item est le résultat d'un processus aléatoire. Par exemple, l'échec enregistré à l'item n° 6 aurait pu, étant donné le score latent qu'on prête à l'enfant (tout en l'ignorant), être une réussite. Autrement dit, on ignore la trame causale de l'échec à cet item à cette date. Le psychologue pourrait suggérer à l'enfant de recommencer la sixième tâche afin que tous deux puissent tenter de comprendre ce qui n'a pas fonctionné.

Quoiqu'il en soit, affirmer que l'enfant est incapable de réussir cet item ne serait pas une proposition valide sans prémisses de laquelle la déduire. Une prémisses convenable serait l'absence d'une condition nécessaire à la réussite de l'item, mais aucun manuel de test ne propose une liste de conditions nécessaires à la réussite des items. La performance au test n'a donc aucune signification diagnostique ni pronostique valide pour cet enfant (Vautier, 2012).

Dans une perspective idiographique, le psychologue tient une position de détective. Il doit identifier les significations que les tâches proposées revêtent pour l'enfant. En particulier, il faut identifier l'intérêt de l'enfant pour les performances tant scolaires que psychotechniques, même si ce type d'interprétation pose de sérieux problèmes méthodologiques. Un enfant dont on peut dire qu'il ne s'engage pas dans la tâche est un enfant qui refuse de montrer ce qu'il sait faire comme ce qu'il ne sait pas faire, auquel cas le psychologue devrait rapporter qu'il n'a pas pu « observer

l'intelligence » de cet enfant parce qu'il ne peut pas assurer que l'enfant s'est « approprié la situation d'examen »<sup>13</sup>.

L'évaluation des aptitudes possède une dimension signifiante complexe. La notation d'une performance qui n'est pas naturellement ordinale constitue une forme d'autorité sociale vis-à-vis de la personne évaluée. Le clinicien doit alors déterminer la liberté dont il jouit, en tant qu'agent de l'évaluation, vis-à-vis de l'impératif évaluatif.

## 5. Conclusion

Les tests d'aptitude ne sont pas des instruments de mesure parce que si c'était vrai, on saurait quels principes ou lois de mesurage opèrent chez les personnes soumises aux items de tests. Ce sont des méthodes d'observation et de scorage de performances. En confondant scorage et mesurage, les psychologues qui défendent la pratique du testage (de l'anglais *testing*) psychologique ne réussiront pas à faire admettre que leurs pratiques sont scientifiquement fondées s'ils s'adressent à des scientifiques, pas plus qu'à dissiper les doutes du profane qui, tout en faisant confiance au professionnalisme des psychologues, ne les crédite pas, à juste titre, d'une science du mesurage. C'est pourquoi la communauté des testeurs n'a pas grand-chose à perdre à reconnaître que l'évaluation psychotechnique soit une sociotechnique, c'est-à-dire un art de préparer des jugements évaluatifs selon les contextes dans lesquels on fait appel à leurs compétences, et n'a aucun devoir de défendre une conception exceptionnelle de ce en quoi consiste le mesurage en psychologie (Michell, 1997, 2000).

En revanche, une telle lucidité dans la communauté des psychologues de tous bords libèrerait la psychologie scientifique de l'impératif de fonder sinon tous, du moins une grande partie de ses programmes de recherche sur la mythologie des grandeurs psychologiques. Si la psychométrie a besoin d'une telle mythologie pour tirer profit de ses modèles à variables latentes, la recherche scientifique en psychologie n'a pas nécessairement besoin des modèles psychométriques (Vautier et al., 2012). Elle a essentiellement besoin du droit à rendre compte de l'immensité de notre ignorance de ce qui cause les comportements qu'on sait observer (Vautier, 2012), laquelle ne peut qu'éclairer la spécificité de la tâche qui consiste à évaluer autrui dans un contexte social quelconque.

## Déclaration d'intérêts

L'auteur déclare ne pas avoir de conflits d'intérêts en relation avec cet article.

## Remerciements

Je remercie Philippe Chartier, Jean-Philippe Gaudron et Valérie Tartas pour leurs commentaires pendant l'élaboration du présent article.

<sup>13</sup> « R6 : L'enfant doit exprimer son accord et s'approprier la situation d'examen » (Voyazopoulos et al., 2011, p. 44). Il semble que ce type de critère nécessite un important travail théorique pour qu'il devienne opérationnel dans une perspective de description « cotation-objective ». Comment fait-on pour décider que l'enfant simulé dans le dialogue « s'approprie la situation d'examen » ? Un psychologue désireux de satisfaire des demandeurs peu sensibles aux subtilités de l'activité mentale, sera plus enclin à juger que l'enfant s'est approprié la situation d'examen qu'un psychologue qui travaille pour des demandeurs moins directifs.



## Références

- Bachelard, G. (1983). *La formation de l'esprit scientifique* (12e ed.). Paris: Vrin.
- Bertrand, D., El Ahmadi, A., & Heuchenne, C. (2008). D'une échelle ordinale de Guttman à une échelle de rapports de Rasch. *Mathématiques et Sciences Humaines*, 4, 25–46.
- Binet, A., & Simon, T. (1907). Le développement de l'intelligence chez les enfants. *L'Année Psychologique*, 14, 1–94.
- Borsboom, D. (2008). Latent variable theory. *Measurement: Interdisciplinary Research and Perspectives*, 6, 25–53.
- Borsboom, D., Cramer, A., Kievit, R. A., Scholten, Z., & Franic, S. (2009). The end of construct validity. In R. W. Lissitz (Ed.), *The concept of validity: Revisions, new directions and applications* (pp. 135–170). Charlotte, NC: Information Age Publishing.
- Canguilhem, G. (1958). Qu'est-ce que la psychologie ? *Revue de Métaphysique et de Morale*, 1, 12–25.
- Chartier, D., & Loarer, E. (2008). *Évaluer l'intelligence logique : Approche cognitive et dynamique*. Paris: Dunod.
- Cizek, G. J. (2012). Defining and distinguishing validity: Interpretations of score meaning and justification of test use. *Psychological Methods*, 17, 31–43.
- Coombs, C. H. (1964). *A theory of data*. New York: Wiley.
- Danziger, K. (1987). Statistical method and the historical development of research practice in American psychology. In L. Krüger, G. Gigerenzer, & M. S. Morgan (Eds.), *Ideas in the sciences* (Vol. 2) *The probabilistic revolution* (pp. 35–47). Cambridge: MIT Press.
- Danziger, K. (1990). *Constructing the subject: Historical origins of psychological research*. New York: Cambridge University Press.
- Duhem, P. (2007). *La théorie physique, son objet, sa structure*. Paris: Vrin.
- Falmagne, J. C. (2003). *Lectures in elementary probability theory and stochastic processes*. Boston: McGraw Hill.
- Fischer, G. H. (1995). Derivations of the Rasch model. In G. H. Fischer, & I. W. Molenaar (Eds.), *Rasch models: foundations, recent developments, and applications* (pp. 15–38). New York: Springer-Verlag.
- Gaillard, F., Colasse, M., Guihard, C., & Michel, R. (2011). Pertinence et nécessité de l'examen psychologique de l'enfant et de l'adolescent. In R. Voyazopoulos, L. Vannetzel, & L.-A. Eynard (Eds.), *L'examen psychologique de l'enfant et utilisation des mesures : Conférence de consensus* (pp. 125–178). Paris: Dunod.
- Granger, G.-G. (1995). *La science et les sciences* (2<sup>e</sup> ed.). Paris: Presses Universitaires de France.
- Grégoire, J. (2009). *L'examen clinique de l'intelligence de l'enfant : fondements et pratique du WISC-IV* (2<sup>e</sup> ed.). Collines de Wavre: Pierre Mardaga Éditeur.
- Guttman, L. (1944). A basis for scaling qualitative data. *American Sociological Review*, 9, 139–150.
- Hacking, I. (2002). *L'émergence de la probabilité*. Paris: Seuil.
- Huteau, M., & Lautrey, J. (1999). Évaluer l'intelligence. In *Psychométrie cognitive*. Paris: Presses Universitaires de France.
- Juhel, J., Gilles, P.-Y., Bouvard, M., Boy, T., Fouques, D., Guimard, P., et al. (2011). Validité des modèles et des outils de l'examen psychologique. In R. Voyazopoulos, L. Vannetzel, & L.-A. Eynard (Eds.), *L'examen psychologique de l'enfant et utilisation des mesures : Conférence de consensus* (pp. 179–251). Paris: Dunod.
- Jumel, B., & Savorinin, F. (2013). *L'aide-mémoire du WISC-IV* (2<sup>e</sup> ed.). Paris: Dunod.
- Kane, M. T. (2006). Validation. In R. L. Brennan (Ed.), *Educational Measurement* (pp. 17–64). Washington, DC: American Council on Education/Praeger.
- Lacot, E., Afzali, M.H., & Vautier, S. (à paraître). Test validation without measurement: Disentangling scientific explanation of item responses and justification of focused assessment policies based on test data. *European Journal of Psychological Assessment*, doi:10.1027/1015-5759/a000253.
- Lamiell, J. T. (1998). 'Nomothetic' and 'idiographic': Contrasting Windelband's understanding with contemporary usage. *Theory & Psychology*, 8, 23–38.
- Michell, J. (1997). Quantitative science and the definition of measurement in psychology. *British Journal of Psychology*, 88, 355–383.
- Michell, J. (2000). Normal science, pathological science and psychometrics. *Theory & Psychology*, 10, 639–667.
- Michell, J. (2003). The quantitative imperative: Positivism, naïve realism, and the place of qualitative methods in psychology. *Theory & Psychology*, 13, 5–31.
- Michell, J. (2009). Invalidity in validity. In R. W. Lissitz (Ed.), *The concept of validity: Revisions, new directions and applications* (pp. 111–133). Charlotte, NC: Information Age Publishing.
- Michell, J. (2013). Constructs, inferences, and mental measurement. *New Ideas in Psychology*, 31, 13–21.
- Newton, P. E. (2012). Clarifying the consensus definition of validity. *Measurement: Interdisciplinary Research and Perspectives*, 10, 1–29.
- Pestre, D. (2013). *À contre-science: Politiques et savoirs des sociétés contemporaines*. Paris: Seuil.

- Popper, K. R. (1973). *La logique de la découverte scientifique*. Paris: Payot.
- Reuchlin, M. (1969). *Les méthodes en psychologie*. Paris: Presses Universitaires de France.
- Salvatore, S., & Valsiner, J. (2010). Between the general and the unique: Overcoming the nomothetic versus idiographic opposition. *Theory & Psychology, 20*, 817–833.
- Searle, J. R. (1995). *The construction of social reality*. New York: The Free Press.
- Vautier, S. (2011). The operationalization of general hypotheses versus the discovery of empirical laws in Psychology. *Philosophia Scientiae, 15*, 105–122.
- Vautier, S. (2012). Propos sur la responsabilité scientifique du psychologue. Essai d'épistémologie appliquée. *Pratiques Psychologiques, 18*, 373–383.
- Vautier, S. (2013). How to state general qualitative facts in psychology? *Quality and Quantity, 47*, 49–56.
- Vautier, S. <http://epistemo.hypotheses.org/1054>
- Vautier, S. <http://epistemo.hypotheses.org/1575>
- Vautier, S. <http://epistemo.hypotheses.org/1608>
- Vautier, S. <http://epistemo.hypotheses.org/520>
- Vautier, S., Lacot, E., & Veldhuis, M. (2014). Puzzle-solving in psychology: The neo-Galtonian vs. nomothetic research focuses. *New Ideas in Psychology, 33*, 46–53.
- Vautier, S., Veldhuis, M., Lacot, E., & Matton, N. (2012). The ambiguous utility of psychometrics for the interpretative founding of socially relevant avatars. *Theory & Psychology, 22*, 810–822.
- Voyazopoulos, R., Vannetzel, L., & Eynard, L.-A. (2011). *L'examen psychologique de l'enfant et l'utilisation des mesures : Conférence de consensus*. Paris: Dunod.
- Wechsler, D. (2005a). *WISC-IV : Manuel d'administration et de cotation*. Paris: Les Éditions du Centre de Psychologie Appliquée.
- Wechsler, D. (2005b). *WISC-IV : Manuel d'interprétation*. Paris: Les Éditions du Centre de Psychologie Appliquée.